# SGML Documents:
# Where Does Quality Go?

José Carlos Ramalho

Jorge Gustavo Rocha

José João Almeida

Pedro Rangel Henriques

Language Processing and Specification Group
Computer Science Department
University of Minho
Portugal

# What will we discuss?

- **When information increases, when information sources increase and vary, what happens to quality?**

- **How can we ensure/preserve quality?**

- **What is quality (what are we talking about)?**

- **In what contexts is quality more relevant?**

- **Can we measure it?          ...**

# What are we doing with SGML?

- **Constructing document DBs**
- **Publishing books on the Internet**
- **Converting parish registers (XIII and XIV century) to SGML**
- **Publishing from SGML DBs: Internet, CDROM, paper, …**
- **Connecting SGML Documents to GIS**

# Quality?

- Quality is good.
- Quality is important.
- Quality is when something is good and achieves to remain good for a period of time.
- Attribute, class, category (from dic.).
- Specific attribute that distinguishes a person, a thing or an entity (from encycolpedia).

# Quality (in our context)?
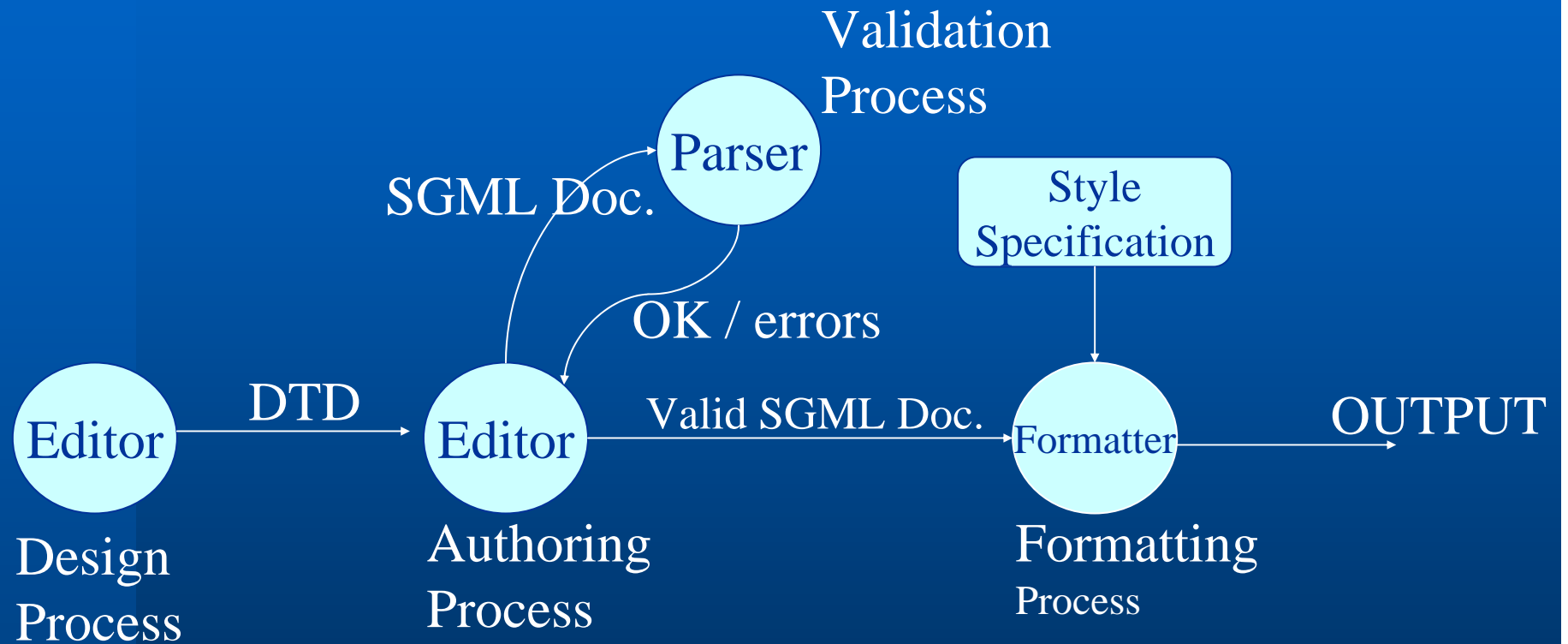
→ **Interface**

→ **…**

→ **Data relevance**

→ **…**

→ **Data correctness**  ⟸  **There is a lot less subjectivity in this item**
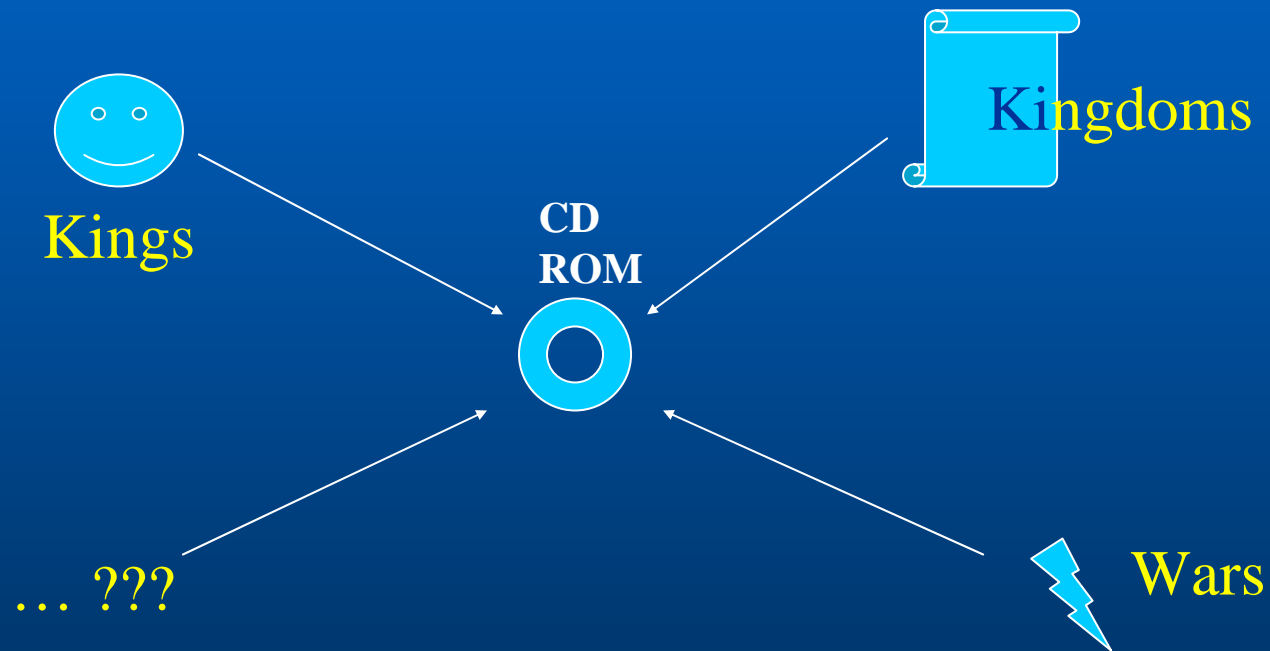
# Aims of this work

- **We want to minimize *Data Incorrectness***
- **We don't want to change existing models**
- **We want to extend them**
- **In the end we want to eliminate information revision cycles**

# SGML authoring and processing model

Validation
Process

Parser

SGML Doc.

Style
Specification

OK / errors

DTD

Editor

Valid SGML Doc.

Editor

Formatter

OUTPUT

Design
Process

Authoring
Process

Formatting
Process

# Data (in)correctness

**Example 1: Portuguese History**

Kings

Kingdoms

CD ROM

… ???

Wars

# Data (in)correctness

**Example 1: Portuguese History**
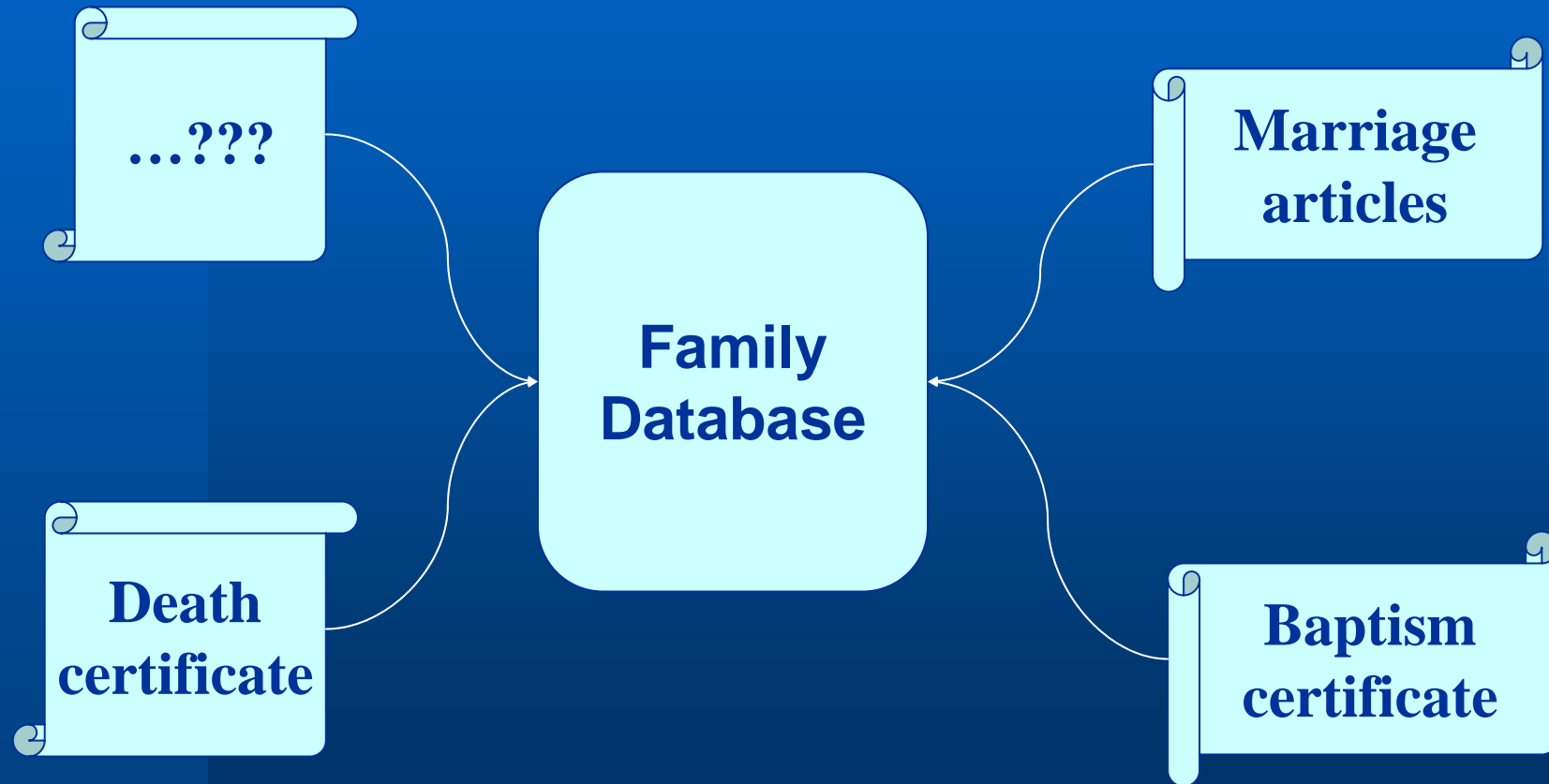
Kingdoms

Kings

CD

Wars

*What went wrong?*

• Kings with inexistent kingdoms

• Wars happening in the wrong era

• Characters that died before they were born

• ...

# Data (in)correctness

Example 2: Parish register (XIII and XIV century)

...???

Marriage articles

Family Database

Death certificate

Baptism certificate

# Data (in)correctness

Birth certificate

Death certificate

**D**

*Problems:*

• negative ages

• death before baptism

• marriages between people with age differences higher than 100

• ...

# What do we propose?

- **An extra validation task:**
  - **we need an additional level of abstraction separating information content from document structure.**
- **Implemented over an external functional system (in the moment …)**
- **Capable of expressing invariants and pre-conditions over data contents**
- **Invisible from the user point of view**

# How?

- **Special** *Comment Sections***: embedding code**
  **in DTDs**

  ```
  <!DOCTYPE  king  [
  <!ELEMENT  king  -- (name,coname, bdate,…)>
  <!-- INV
     inv_king(k) = …
  -->
  ```

- **Throught an anchor to an external file**

  ```
  <!-- INV: king.cam -->
  <!DOCTYPE  king [ … ]>
  ```

# Example: kings and decrees

```
<!-- INV: king.cam  -->
<!DOCTYPE  king  [
<!ELEMENT  king -- (name, coname,
    bdate, ddate,decree+)>
<!ELEMENT  decree -- (date, body)>
<!ELEMENT
    (name,coname,bdate,ddate,date) --
    (#PCDATA)>
<!ELEMENT  body -- (#PCDATA)>
]>
```

king.dtd

$$Inv\_king(k) =$$
$$\{ \ if( \ k \ notin \ famous\_personsDB \rightarrow$$
$$k \ ++ \ `` \ not \ in \ FPDB''),$$
$$if( \ bdate\_(k) > ddate\_(k) \rightarrow k \ ++$$
$$\text{``died before he has born''}),$$
$$if( \ ddate\_(k) - bdate\_(k) > 120 \rightarrow$$
$$k \ ++ \ \text{``lived more than 120''}),$$
$$if( \ !all( \ x \leftarrow decree\_l(k) :$$
$$bdate\_(k) < date\_(x) \ \wedge$$
$$date\_(x) < ddate\_(k) \ ) \rightarrow$$
$$k \ ++ \ \text{``made a decree outside}$$
$$\text{his life'' } )$$
$$\};$$

king.cam

# Example: kings and decrees

```
<king>
    <name>D.Dinis</name>
    <coname>Farmer</coname>
    <bdate>1270.09.23</bdate>
    <ddate>1370.09.23</ddate>
    <decree>
        <date>1300.07.15</date>
        <body>From this day only
 bicycles are allowed to
 circulate.</body>
    </decree>
    <decree>
        <date>1389.11.03</date>
        <body>McDonald's will sell
 green wine instead of COCA-
COLA.</decree>
</king>
```
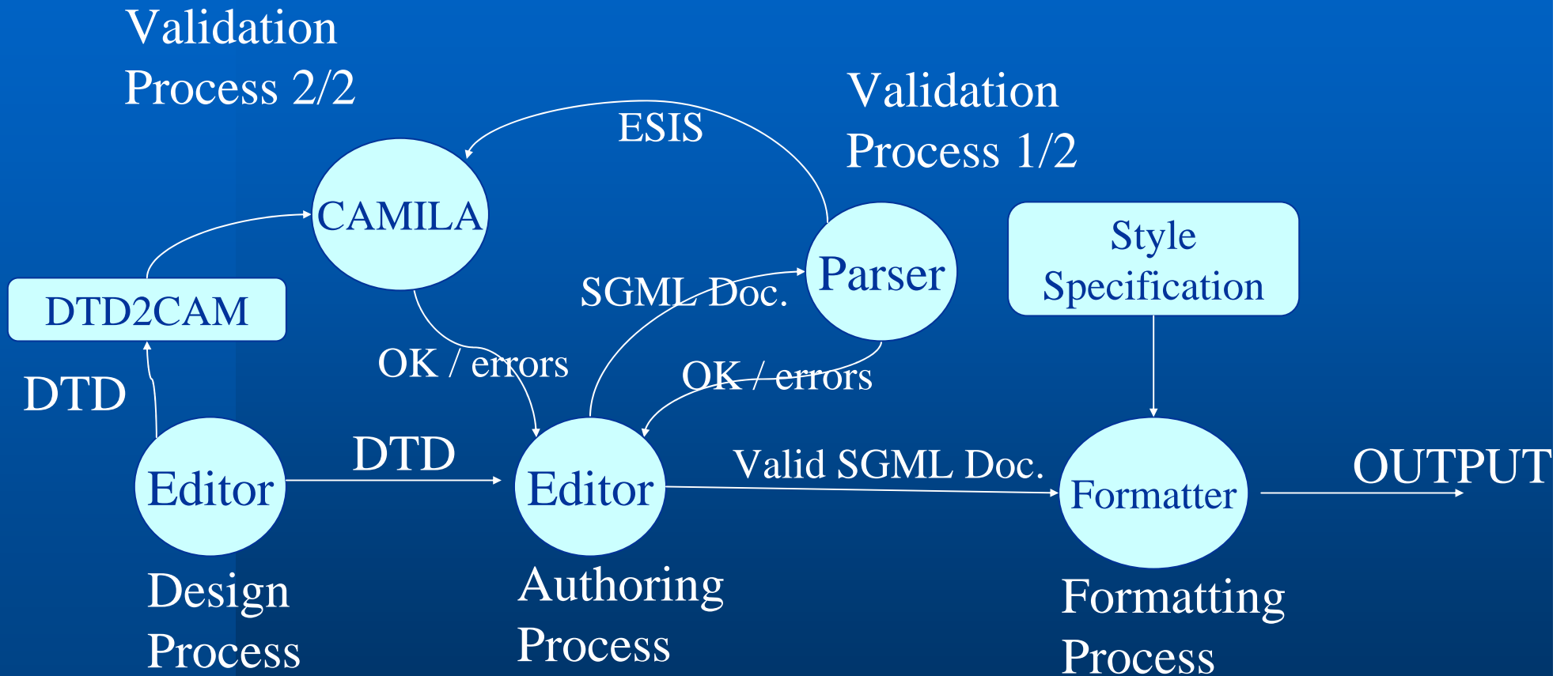
ERRORS:

D.Dinis  must be inserted in FPDB.

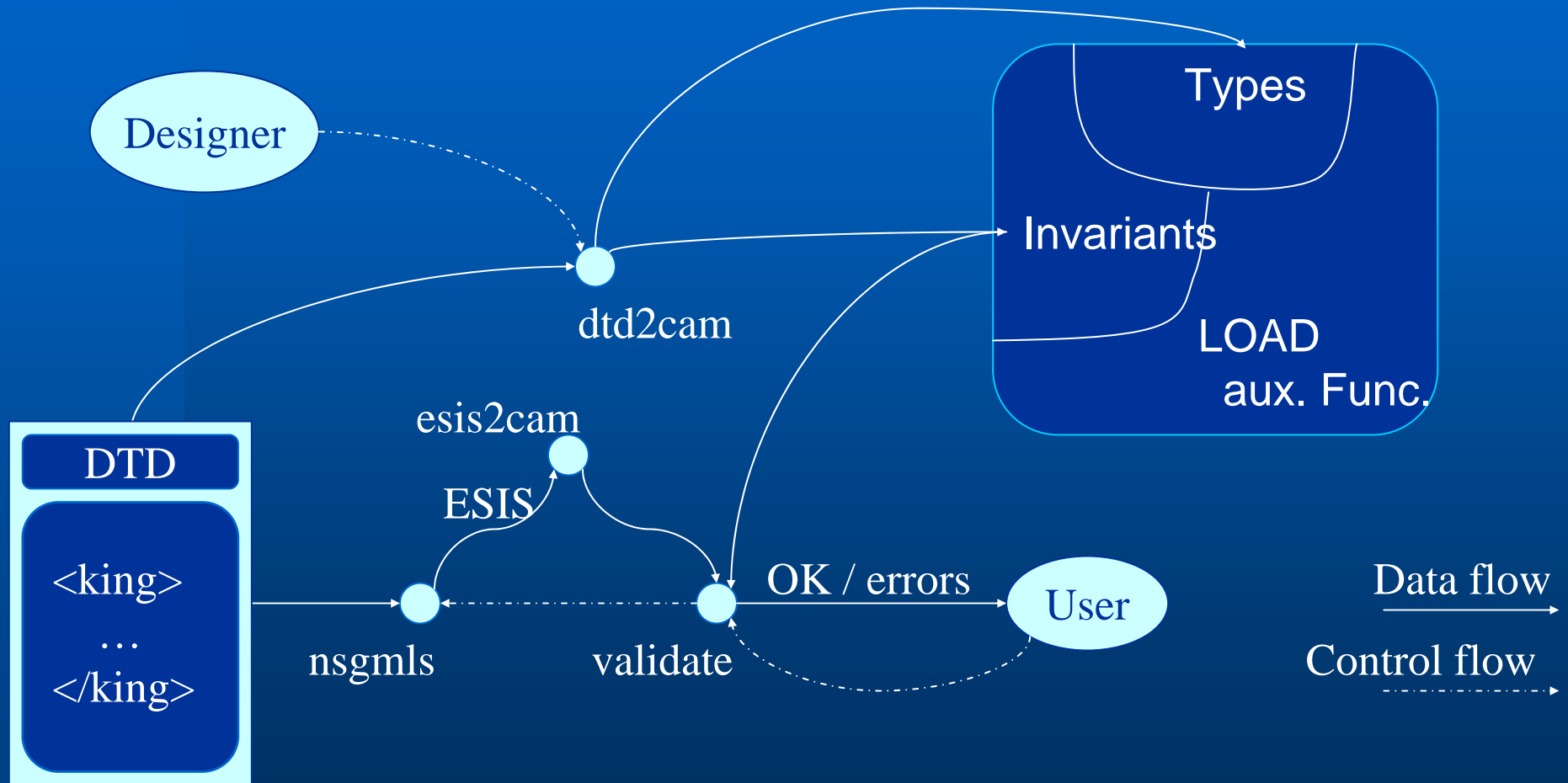D.Dinis made a decree outside his life.

# Other Examples

- **Tying an Archaeological Database to a GIS:**
  - archaeological SGML documents have geographical coordinates.
  - we must ensure that every one of those coordinates is within a certain range.

- **City Council Elections**
  - each voting section produces a final report with the results (an SGML document).
  - we must ensure that the number of votes matches the number of subscribed voters minus the absent ones.

# *New* SGML auth. and proc. model

Validation
Process 2/2

Validation
Process 1/2

ESIS

CAMILA

Style
Specification

DTD2CAM

SGML Doc.

Parser

DTD

OK / errors

OK / errors

Editor

DTD

Editor

Valid SGML Doc.

Formatter

OUTPUT

Design
Process

Authoring
Process

Formatting
Process

# *Camila* Validation Process



Designer

Types

Invariants

LOAD
aux. Func.

dtd2cam

esis2cam

DTD

ESIS

<king>
…
</king>

nsgmls

validate

OK / errors

User

Data flow

Control flow

# *Camila* Validation Process

Designer

dtd2cam

Types

Invariants

LOAD

<!ELEMENT  king  - - (name, coname,

bdate, ddate, decree+)>

dtd2cam

TYPE
   king = name_ :name
         coname_ :coname
         bdate_ :bdate
         ddate_ :ddate
         decree_l :decree-seq
;
ENDTYPE
inv_king( k ) = true;

Data flow

rol flow

# Conclusion

- **The new proposed model enables us to put some kind of data constraints associated with DTD element contents.**

- **We can avoid many errors given by a distracted user.**

- **We can improve information quality and reduce information revision cycle.**
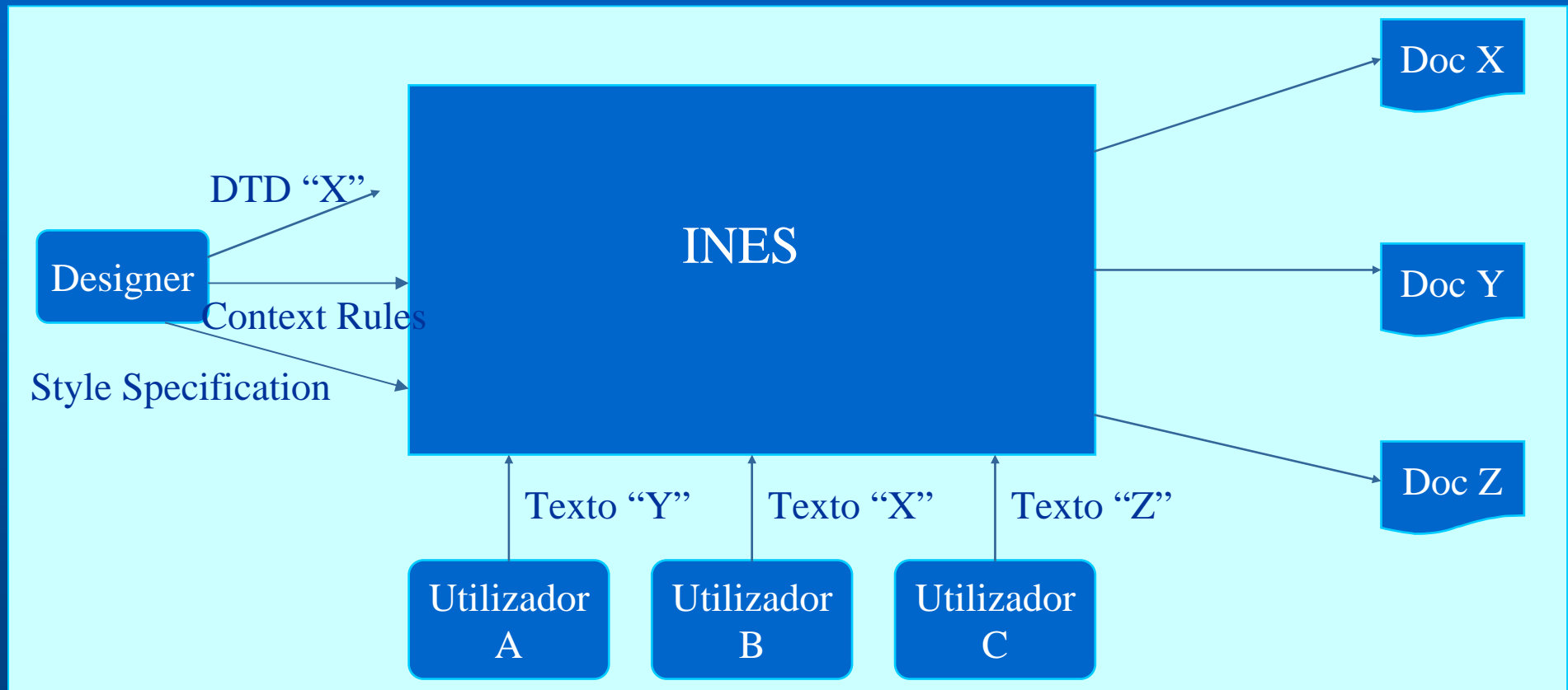
# Conclusion (cont.)

- **In the case studies we have dealed with so far we didn't find complex invariants.**

- **Structural correctness imposed by SGML already enforces some validation over element contents.**

- **Most of needed invariants are very simple: domain range validation, relationship validation, ...**

# Future Work

- **A simple constraint language is being studied/created to optimize the proposed system.**

- **We are going to implement this validation scheme (with the new language) in our prototype INES ("A Document Programming Environment").**

# INES: Document Programming Env

# INES: inside

SGML text

Errors

Designer

Context Conditions;
Invariants

Style
Specification

DTD
Editor

Código Scheme

DTD

Editor
Generator

SGEN

RTF

PostScript

DTD

Context
Editor

"X"
Editor

Doc X

DSSSL
Editor

Text

Errors

Utilizador